

Title	確率伝搬法による近似と情報幾何学による解析(情報物理学の数学的構造)
Author(s)	Ikeda, Shiro
Citation	数理解析研究所講究録 (2007), 1532: 1-10
Issue Date	2007-02
URL	http://hdl.handle.net/2433/58962
Right	
Type	Departmental Bulletin Paper
Textversion	publisher

確率伝搬法による近似と情報幾何学による解析

統計数理研究所 池田思朗 (Shiro Ikeda)
The Institute of Statistical Mathematics

1 Introduction

We have developed an information geometrical framework of the Belief propagation (BP) algorithm in [1, 2]. One important aspect of the information geometrical framework is that we can view different methods developed in different fields from a unified viewpoint. In this draft, we first summarize the information geometrical framework, and show that the adaptive TAP[3] method is also expressed in the same framework. This will allow us to apply well-developed tools and results of the loopy BP to the adaptive TAP.

2 Information geometrical framework of BP

First, we summarize the information geometrical view of BP.

2.1 Problem

The distribution we would like to focus on is given as follows

$$q(\mathbf{x}) = \exp[c_0(\mathbf{x}) + c_1(\mathbf{x}) + \cdots + c_L(\mathbf{x}) - \psi_q].$$

For a while, we restrict ourselves to the case where $\mathbf{x} \in \{-1, +1\}^N$. In the case of the Boltzmann machine,

$$c_0(\mathbf{x}) = \mathbf{h} \cdot \mathbf{x}, \quad c_r(\mathbf{x}) = J_{ij} x_i x_j.$$

For a large N , we cannot compute ψ_q , which is the log-partition function.

The goal of the BP algorithm is to infer $\boldsymbol{\eta} = E_q[\mathbf{x}]$, or equivalently compute $\prod_i q(x_i)$. The direct computation is not tractable when N is large and the graph is cyclic.

2.2 Models, manifolds, and equilibrium

Let us define the set of all distributions S .

$$S = \left\{ p(\mathbf{x}) \mid \sum_{\mathbf{x}} p(\mathbf{x}) = 1, p(\mathbf{x}) > 0 \right\}.$$

A submanifold $M_0 \subset S$ is defined as

$$M_0 = \left\{ p_0(\mathbf{x}; \boldsymbol{\theta}) = \exp[\mathbf{h} \cdot \mathbf{x} + \boldsymbol{\theta} \cdot \mathbf{x} - \psi_0(\boldsymbol{\theta})] \mid \mathbf{h}, \boldsymbol{\theta} \in \mathbb{R}^N \right\},$$

where \cdot shows the inner product. This is a set of factorizable distributions. Each component is independent for the distributions of M_0 , and its natural parameter is $\boldsymbol{\theta}$. Conversely, every factorizable distribution is included in M_0 . Therefore, the problem to compute $\prod_i q(x_i)$ is equivalent to compute the natural parameter $\boldsymbol{\theta}$ which satisfies $p_0(\mathbf{x}; \boldsymbol{\theta}) = \prod_i q(x_i)$.

Let us define the m -projection of distribution $r(\mathbf{x})$ to M_0 as

$$\boldsymbol{\theta} = \pi_{M_0} \circ r(\mathbf{x}) = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} D[r(\mathbf{x}); p_0(\mathbf{x}; \boldsymbol{\theta})], \quad \Pi_{M_0} \circ r(\mathbf{x}) = \underset{p(\mathbf{x}) \in M_0}{\operatorname{argmin}} D[r(\mathbf{x}); p_0(\mathbf{x}; \boldsymbol{\theta})],$$

where $D[\cdot; \cdot]$ is the KL-divergence, and we have the following relation

$$p_0(\mathbf{x}; \boldsymbol{\theta}^*) = \Pi_{M_0} \circ q(\mathbf{x}) = \prod_{i=1}^N q(x_i), \quad \text{where } \boldsymbol{\theta}^* = \pi_{M_0} \circ q(\mathbf{x}).$$

Now, let us define $p_r(\mathbf{x}; \boldsymbol{\zeta}_r)$, $r = 1, \dots, L$ as,

$$p_r(\mathbf{x}; \boldsymbol{\zeta}_r) = \exp[\mathbf{h} \cdot \mathbf{x} + c_r(\mathbf{x}) + \boldsymbol{\zeta}_r \cdot \mathbf{x} - \psi_r(\boldsymbol{\zeta}_r)], \quad \boldsymbol{\zeta}_r \in \mathbb{R}^N, \quad r = 1, \dots, L.$$

$p_r(\mathbf{x}; \boldsymbol{\zeta}_r)$ is an exponential family which includes $c_r(\mathbf{x})$.

$$M_r = \left\{ p_r(\mathbf{x}; \boldsymbol{\zeta}_r) \mid \boldsymbol{\zeta}_r \in \mathbb{R}^N \right\}, \quad r = 1, \dots, L.$$

Its natural parameter is $\boldsymbol{\zeta}_r$. We assume the computation of $\pi_{M_0} \circ p_r(\mathbf{x}; \boldsymbol{\zeta}_r)$ is tractable for every $\boldsymbol{\zeta}_r \in \mathbb{R}^N$ and r . With $p_0(\mathbf{x}; \boldsymbol{\theta})$ and $p_r(\mathbf{x}; \boldsymbol{\zeta}_r)$, $r = 1, \dots, L$, the BP algorithm is defined as follows,

1. Set $t = 0$, $\boldsymbol{\xi}_r^t = \mathbf{0}$, $\boldsymbol{\zeta}_r^t = \mathbf{0}$, $r = 1, \dots, L$.

2. Increase t by 1 and update $\boldsymbol{\xi}_r^{t+1}$, $r = 1, \dots, L$ as follows

$$\boldsymbol{\xi}_r^{t+1} = \pi_{M_0} \circ p_r(\mathbf{x}; \boldsymbol{\zeta}_r^t) - \boldsymbol{\zeta}_r^t.$$

3. Update $\boldsymbol{\theta}^{t+1}$ and $\boldsymbol{\zeta}_r^{t+1}$ as follows

$$\boldsymbol{\zeta}_r^{t+1} = \sum_{r' \neq r} \boldsymbol{\xi}_{r'}^{t+1}, \quad \boldsymbol{\theta}^{t+1} = \sum_r \boldsymbol{\xi}_r^{t+1} = \frac{1}{L-1} \sum_r \boldsymbol{\zeta}_r^{t+1}.$$

4. Repeat 2 and 3 until $\{\boldsymbol{\xi}_r^t\}$ converges.

At the equilibrium, we know the following conditions hold [2].

m -condition: $\theta^* = \pi_{M_0} \circ p_r(\mathbf{x}; \zeta_r^*)$.

e -condition: $\theta^* = \frac{1}{L-1} \sum_{r=1}^L \zeta_r^*$, or equivalently $q(\mathbf{x}) = \frac{1}{Z} \frac{\prod_{r=1}^L p_r(\mathbf{x}; \zeta_r^*)}{p_0(\mathbf{x}; \theta^*)^{L-1}}$.

We have shown these conditions are satisfied at the equilibrium of BP, TRP, and CCCP in [2].

2.3 Free energy

When the e -condition is satisfied, we have the following relation.

$$q(\mathbf{x}) = \frac{1}{Z} \frac{\prod_{r=1}^L p_r(\mathbf{x}; \zeta_r^*)}{p_0(\mathbf{x}; \theta^*)^{L-1}}.$$

Moreover, when graph is tree, at the equilibrium of BP, Z becomes 1, that is

$$q(\mathbf{x}) = \frac{\prod_{r=1}^L p_r(\mathbf{x}; \zeta_r^*)}{p_0(\mathbf{x}; \theta^*)^{L-1}}.$$

This equation shows that when graph is tree, free energy is expressed as

$$\mathcal{F} = -\psi_q = (L-1)\psi_0(\theta^*) - \sum_{r=1}^L \psi_r(\zeta_r^*).$$

This relation does not hold for cyclic graphs, but the right hand side of the equation is called the Bethe free energy. This equation is deeply related to the BP algorithm [2].

$$\mathcal{F}_{\text{Bethe}} = (L-1)\psi_0(\theta^*) - \sum_{r=1}^L \psi_r(\zeta_r^*).$$

3 Adaptive TAP

The adaptive TAP method is described in detail by Oppor and Winther[3]. We follow their results, where notations are slightly modified to make consistency with this memo.

3.1 Problem

We consider the distribution

$$q(\mathbf{x}) = \frac{1}{Z} \prod_{r=1}^N \rho(x_r) \exp \left[\mathbf{h} \cdot \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{J} \mathbf{x} \right]. \quad (1)$$

\mathbf{J} is a symmetric matrix where diagonal elements are 0. In the paper, ρ takes a lot of kinds of functions, such as Dirac delta function. Actually, by taking $\rho(x_r) = (\delta(x_r - 1) + \delta(x_r + 1))/2$, $q(\mathbf{x})$ becomes equivalent to the Boltzmann machine.

3.2 Adaptive TAP equations

The aim of the adaptive TAP approach is to infer $E_q[x_r]$ and $E_q[x_r^2]$. Let m_r be the inference of $E_q[x_r]$. A summary of Adaptive TAP equations is given as follow.

$$m_r = \frac{\partial}{\partial h_r} \ln Z_0^{(r)} \quad (2)$$

$$Z_0^{(r)} = \int \rho(x_r) \exp \left[\left(\sum_{s=1}^N J_{rs} m_s - V_r m_r + h_r \right) x_r + \frac{V_r}{2} x_r^2 \right] dx_r$$

$$\frac{\partial m_r}{\partial h_r} = \frac{\partial^2}{\partial h_r^2} \ln Z_0^{(r)} = [(S - J)^{-1}]_{rr}, \quad (3)$$

$$S = \text{diag}(s_1, \dots, s_N), \quad \frac{1}{s_r - V_r} = \frac{\partial m_r}{\partial h_r} = [(S - J)^{-1}]_{rr}.$$

Let us define $p_r(x_r)$ as follows,

$$p_r(x_r) = \frac{1}{Z_0^{(r)}} \rho(x_r) \exp \left[\left(\sum_{s=1}^N J_{rs} m_s - V_r m_r + h_r \right) x_r + \frac{V_r}{2} x_r^2 \right]. \quad (4)$$

From the definition of $Z_0^{(r)}$, $p_r(x_r)$ is a density function of x_r . Now, we can rewrite the adaptive TAP equations (2) and (3) as follows,

$$m_r = \int x_r p_r(x_r) dx_r \quad (5)$$

$$[(S - J)^{-1}]_{rr} = \int (x_r - m_r)^2 p_r(x_r) dx_r = \int x_r^2 p_r(x_r) dx_r - m_r^2. \quad (6)$$

3.3 Free energy

The free energy, which is the inference of $-\ln Z$ is given as follows at the solutions of adaptive TAP equations.

$$\begin{aligned} \Phi_1 &= \Phi_0 - \frac{1}{2} \mathbf{m}^T J \mathbf{m} + \frac{1}{2} \ln \det(S - J) - \frac{1}{2} \sum_{r=1}^N V_r \chi_{rr} + \frac{1}{2} \sum_{r=1}^N \ln \chi_{rr} \\ \Phi_0 &= - \sum_{r=1}^N \ln Z_0^{(r)} + \mathbf{m}^T J \mathbf{m} + \frac{1}{2} \sum_{r=1}^N V_r M_r, \\ \chi_{rr} &= M_r - m_r^2 = [(S - J)^{-1}]_{rr}. \end{aligned}$$

where M_r is the inference of $E_q[x_r^2]$. We can simplify the above free energy as

$$\Phi_1 = - \sum_{r=1}^N \ln Z_0^{(r)} + \frac{1}{2} \ln \det(S - J) + \frac{1}{2} \sum_{r=1}^N \ln [(S - J)^{-1}]_{rr} + \frac{1}{2} \mathbf{m}^T J \mathbf{m} - \frac{1}{2} \sum_{r=1}^N V_r m_r^2. \quad (7)$$

4 From BP to adaptive TAP

We show that the information geometrical framework of the BP algorithm explains the adaptive TAP equations and free energy.

4.1 Problem

We consider the case where $\rho(x_r)$ is strictly positive for $x_r \in \mathbb{R}$. This is not true when ρ is the Dirac delta function. Naively, we can treat such a problem by putting

$$\rho(x_r) = \frac{1}{2\sqrt{2\pi\sigma^2}} \left(\exp\left[-\frac{(x_r-1)^2}{2\sigma^2}\right] + \exp\left[-\frac{(x_r+1)^2}{2\sigma^2}\right] \right),$$

and bringing $\sigma^2 \rightarrow 0$. But we put this problem beside and assume $\rho(x_r) > 0$. Equation (1) can be rewritten as

$$\begin{aligned} q(\mathbf{x}) &= \exp[c_0(\mathbf{x}) + c_1(\mathbf{x}) + \dots + c_N(\mathbf{x}) - \psi_q], \\ c_0(\mathbf{x}) &= \frac{1}{2}\mathbf{x}^T J \mathbf{x}, \quad c_r(\mathbf{x}) = c_r(x_r) = \ln \rho(x_r) + h_r x_r, \quad \psi_q = \ln Z. \end{aligned} \quad (8)$$

Next, we define p_0 as the distribution whose sufficient statistics are $x_r, x_r^2, r = 1, \dots, N$. Natural choice is the normal distribution, which is defined as

$$\begin{aligned} p_0(\mathbf{x}; \boldsymbol{\mu}, S) &= \exp\left[c_0(\mathbf{x}) + \boldsymbol{\mu} \cdot \mathbf{x} - \frac{1}{2}\mathbf{x}^T S \mathbf{x} - \psi_0(\boldsymbol{\mu}, S)\right], \\ S &= \text{diag}(s_1, \dots, s_N), \\ \psi_0(\boldsymbol{\mu}, S) &= \frac{N}{2} \ln 2\pi - \ln \det(S - J) + \frac{1}{2}\boldsymbol{\mu}^T (S - J)^{-1} \boldsymbol{\mu}. \end{aligned} \quad (9)$$

From the definition $p_0(\mathbf{x}; \boldsymbol{\mu}, S) \sim \mathcal{N}((S - J)^{-1} \boldsymbol{\mu}, (S - J)^{-1})$, where \mathcal{N} shows the density function of a normal distribution. We set M_0 as

$$M_0 = \{p_0(\mathbf{x}; \boldsymbol{\mu}, S) \mid \boldsymbol{\mu} \in \mathbb{R}^N, S = \text{diag}(s_1, \dots, s_N), s_i > 0\}$$

Now, our ultimate goal is to obtain the m -projection of $q(\mathbf{x})$ defined as (8) to M_0 , which then provides the exact mean and variance of x_r .

4.2 e -condition and m -condition

We define $p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r})$ as follows

$$\begin{aligned} p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) &= \exp\left[c_0(\mathbf{x}) + c_r(x_r) + \boldsymbol{\mu}_{\setminus r} \cdot \mathbf{x}_{\setminus r} - \frac{1}{2}\mathbf{x}_{\setminus r}^T S_{\setminus r} \mathbf{x}_{\setminus r} - \psi_r(\boldsymbol{\mu}_{\setminus r}, S_{\setminus r})\right], \\ S_{\setminus r} &= \text{diag}(s_1, \dots, s_{r-1}, s_{r+1}, \dots, s_N) \in \mathbb{R}^{(N-1) \times (N-1)}, \\ \boldsymbol{\mu}_{\setminus r} &= (\mu_1, \dots, \mu_{r-1}, \mu_{r+1}, \dots, \mu_N)^T \in \mathbb{R}^{N-1}, \\ \mathbf{x}_{\setminus r} &= (x_1, \dots, x_{r-1}, x_{r+1}, \dots, x_N)^T \in \mathbb{R}^{N-1}. \end{aligned}$$

We define M_r as

$$M_r = \{p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r})\}.$$

The difference between p_0 and p_r is that,

$$“p_r \text{ does not include } \mu_r x_r \text{ nor } -s_r x_r^2/2, \text{ but includes } c_r(x_r).”$$

This is similar to the information geometrical framework of the BP algorithm.

e-condition

It is easy to show the e-condition holds, that is

$$q(\mathbf{x}) = \frac{1}{Z} \frac{\prod_{r=1}^N p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r})}{p_0(\mathbf{x}; \boldsymbol{\mu}, S)^{N-1}}. \quad (10)$$

The numerator of the right hand side is

$$\prod_{r=1}^N p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) \propto \exp \left[N c_0(\mathbf{x}) + \sum_{r=1}^N c_r(x_r) + (N-1) \sum_{r=1}^N \mu_r x_r - \frac{(N-1)}{2} \sum_{r=1}^N s_r x_r^2 \right],$$

and its denominator is

$$p_0(\mathbf{x}; \boldsymbol{\mu}, S)^{N-1} \propto \exp \left[(N-1) c_0(\mathbf{x}) + (N-1) \sum_{r=1}^N \mu_r x_r - \frac{(N-1)}{2} \sum_{r=1}^N s_r x_r^2 \right],$$

which proves (10).

m-condition

Next, we show that the m -condition corresponds to the adaptive TAP equations. When the m -condition holds,

$$p_0(\mathbf{x}; \boldsymbol{\mu}, S) = \Pi_{M_0} \circ p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) = \underset{p(\mathbf{x}) \in M_0}{\operatorname{argmin}} D[p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}); p(\mathbf{x})], \quad r = 1, \dots, N.$$

The sufficient statistics of p_0 is $x_r, x_r^2, r = 1, \dots, N$. Therefore the m -condition is equivalent to

$$\begin{aligned} \mathbf{m} &= K^{-1} \boldsymbol{\mu} = \int \mathbf{x} p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) d\mathbf{x}, \\ \frac{1}{[K^{-1}]_{ss}} &= \int x_s^2 p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) d\mathbf{x} - m_s^2, \quad r, s = 1, \dots, N. \\ \text{where, } K &= S - J \end{aligned}$$

Let us set \mathbf{J}_r , $K_{\setminus r}$, and $J_{\setminus r}$ as follows

$$J_{\setminus r} = \begin{pmatrix} 0 & \cdots & J_{1(r-1)} & J_{1(r+1)} & \cdots & J_{1N} \\ \vdots & \ddots & & & & \vdots \\ J_{(r-1)1} & & & & & J_{(r-1)N} \\ J_{(r+1)1} & & & & & J_{(r+1)N} \\ \vdots & & & & \ddots & \vdots \\ J_{N1} & \cdots & J_{N(r-1)} & J_{N(r+1)} & \cdots & 0 \end{pmatrix} \in \mathbb{R}^{(N-1) \times (N-1)}$$

$$K_{\setminus r} = (S_{\setminus r} - J_{\setminus r}) \in \mathbb{R}^{(N-1) \times (N-1)}$$

$$\mathbf{J}_r = (J_{r1}, \dots, J_{r(r-1)}, J_{r(r+1)}, \dots, J_{rN})^T \in \mathbb{R}^{N-1}$$

We can rewrite $p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r})$ as follows,

$$\begin{aligned} p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) &= p_r(x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) p_r(\mathbf{x}_{\setminus r} | x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) \\ &= p_r(x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) \mathcal{N}(K_{\setminus r}^{-1}(\boldsymbol{\mu}_{\setminus r} + \mathbf{J}_r x_r), K_{\setminus r}^{-1}) \end{aligned} \quad (11)$$

where,

$$p_r(x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) = \frac{1}{Z_r} \exp \left[c_r(x_r) + \frac{1}{2} (\mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r) x_r^2 + (\mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r}) x_r \right]. \quad (12)$$

We can also rewrite $p_0(\mathbf{x}; \boldsymbol{\mu}, S)$ as follows,

$$\begin{aligned} p_0(\mathbf{x}; \boldsymbol{\mu}, S) &= p_0(x_r; \boldsymbol{\mu}, S) p_0(\mathbf{x}_{\setminus r} | x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) \\ &= p_0(x_r; \boldsymbol{\mu}, S) \mathcal{N}(K_{\setminus r}^{-1}(\boldsymbol{\mu}_{\setminus r} + \mathbf{J}_r x_r), K_{\setminus r}^{-1}) \end{aligned} \quad (13)$$

where,

$$\begin{aligned} p_0(x_r; \boldsymbol{\mu}, S) &= \frac{1}{Z_r} \exp \left[-\frac{1}{2} s_r x_r^2 + \mu_r x_r + \frac{1}{2} (\mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r) x_r^2 + (\mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r}) x_r \right] \\ &= \mathcal{N} \left(\frac{1}{s_r - \mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r} (\mu_r + \mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r}), \frac{1}{s_r - \mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r} \right) \end{aligned}$$

And we have the following proposition

proposition 1. *The m -condition is satisfied, if and only if the following equations hold.*

$$m_r = \int x_r p_r(x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) dx_r \quad (14)$$

$$[(S - J)^{-1}]_{rr} = \frac{1}{s_r - \mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r} = \int x_r^2 p_r(x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) dx_r - m_r^2 \quad (15)$$

$$r = 1, \dots, N, \quad \text{where } \mathbf{m} = (S - J)^{-1} \boldsymbol{\mu}.$$

Now we move to the adaptive TAP equations

lemma 1. *Equations (14) and (15) are equivalent to (5) and (6).*

Proof. What we have to show is that $p_r(x_r)$ in (4) is equivalent to $p_r(x_r; \mu_{\setminus r}, S_{\setminus r})$ in (12). We show both of them again.

$$p_r(x_r) = \frac{1}{Z_0^{(r)}} \rho(x_r) \exp \left[\left(\sum_{j=1}^N J_{rj} m_j - V_r m_r + h_r \right) x_r + \frac{V_r}{2} x_r^2 \right] \quad (16)$$

$$\begin{aligned} p_r(x_r; \mu_{\setminus r}, S_{\setminus r}) &= \frac{1}{Z_r} \exp \left[c_r(x_r) + \frac{1}{2} (\mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r) x_r^2 + (\mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r}) x_r \right] \\ &= \frac{1}{Z_r} \rho(x_r) \exp \left[(\mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r} + h_r) x_r + \frac{1}{2} (\mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r) x_r^2 \right] \end{aligned} \quad (17)$$

$$[(S - J)^{-1}]_{rr} = \frac{1}{s_r - V_r}, \quad \mathbf{m}_{\setminus r} = (m_1, \dots, m_{r-1}, m_{r+1}, \dots, m_N)^T \in \mathfrak{R}^{N-1}. \quad (18)$$

where we used $c_r(x_r) = \ln \rho(x_r) + h_r x_r$. And by comparing (16) and (17), what we have to prove is

$$\mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r} = \sum_{s=1}^N J_{rs} m_s - V_r m_r = \mathbf{J}_r^T \mathbf{m}_{\setminus r} - V_r m_r \quad (19)$$

$$\mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r = V_r. \quad (20)$$

From (15) and (18), $V_r = \mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r$ is shown and (20) is proved, and from the relation $\mathbf{m} = K^{-1} \boldsymbol{\mu}$ or $\boldsymbol{\mu} = K \mathbf{m}$, we have

$$\begin{aligned} \begin{pmatrix} \mu_r \\ \boldsymbol{\mu}_{\setminus r} \end{pmatrix} &= \begin{pmatrix} s_r & -\mathbf{J}_r^T \\ -\mathbf{J}_r & K_{\setminus r} \end{pmatrix} \begin{pmatrix} m_r \\ \mathbf{m}_{\setminus r} \end{pmatrix} \\ \mu_r &= s_r m_r - \mathbf{J}_r^T \mathbf{m}_{\setminus r} \\ \boldsymbol{\mu}_{\setminus r} &= -\mathbf{J}_r m_r + K_{\setminus r} \mathbf{m}_{\setminus r}. \end{aligned}$$

Which proves (19), that is,

$$\begin{aligned} \mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r} &= -\mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r m_r + \mathbf{J}_r^T \mathbf{m}_{\setminus r} \\ &= \mathbf{J}_r^T \mathbf{m}_{\setminus r} - V_r m_r. \end{aligned}$$

Thus, the m -condition derives the adaptive TAP equations. \square

4.3 Free energy

Since we have seen that the framework of the adaptive TAP method is very similar to that of the BP algorithm, we can define the free energy as the Bethe free energy, that is

$$\mathcal{F} = (N - 1) \psi_0(\boldsymbol{\mu}, S) - \sum_{r=1}^N \psi_r(\boldsymbol{\mu}_{\setminus r}, S_{\setminus r}).$$

We now show that this free energy is equivalent to the adaptive TAP free energy in (7). Above equation can be rewritten as

$$\mathcal{F} = \sum_{r=1}^N (\psi_0(\boldsymbol{\mu}, S) - \psi_r(\boldsymbol{\mu}_{\setminus r}, S_{\setminus r})) - \psi_0(\boldsymbol{\mu}, S).$$

From (9),

$$\begin{aligned}\psi_0(\boldsymbol{\mu}, S) &= \frac{N}{2} \ln 2\pi - \ln \det(S - J) + \frac{1}{2} \boldsymbol{\mu}^T (S - J)^{-1} \boldsymbol{\mu} \\ &= \frac{N}{2} \ln 2\pi - \ln \det(S - J) + \frac{1}{2} \mathbf{m}^T (S - J) \mathbf{m}\end{aligned}$$

And by comparing (11) and (13), we have

$$\begin{aligned}\psi_0(\boldsymbol{\mu}, S) - \psi_r(\boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) &= -\ln Z_r + \frac{1}{2} \ln 2\pi - \frac{1}{2} \ln(s_r - \mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r) \\ &\quad + \frac{1}{s_r - \mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r} (\mu_r + \mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r})^2 \\ &= -\ln Z_r + \frac{1}{2} \ln 2\pi + \frac{1}{2} \ln[(S - J)^{-1}]_{rr} + \frac{1}{2} \frac{1}{[(S - J)^{-1}]_{rr}} m_r^2\end{aligned}$$

The free energy is rewritten as

$$\begin{aligned}\mathcal{F} &= -\sum_{r=1}^N \ln Z_r + \ln \det(S - J) - \frac{1}{2} \mathbf{m}^T (S - J) \mathbf{m} + \frac{1}{2} \sum_{r=1}^N \ln[(S - J)^{-1}]_{rr} + \frac{1}{2} \sum_{r=1}^N \frac{1}{[(S - J)^{-1}]_{rr}} m_r^2 \\ &= -\sum_{r=1}^N \ln Z_r + \ln \det(S - J) + \frac{1}{2} \sum_{r=1}^N \ln[(S - J)^{-1}]_{rr} + \frac{1}{2} \mathbf{m}^T J \mathbf{m} - \frac{1}{2} \sum_{r=1}^N V_r m_r^2.\end{aligned}$$

This is equivalent to (7).

4.4 BP like algorithm for adaptive TAP

It is written in [3], that it is not easy to compute the solution of adaptive TAP equations. But since the framework of adaptive TAP method is similar to that of the BP, we can derive BP or related algorithms such as CCCP to solve the adaptive TAP equations. Here we show one example, which is the BP like algorithm to solve the adaptive TAP equations.

1. Set $t = 0$ and $\boldsymbol{\mu}^{(t)} = \mathbf{o}$, $s_r^{(t)} = c$, ($c > 0$, for example, $c = 1$) $r = 1, \dots, L$.
2. Increase t by 1 and calculate $m_r^{(t)}$ and $\sigma_r^{(t)2}$ $r = 1, \dots, L$ as follows

$$\begin{aligned}m_r^{(t)} &= \int x_r p_r(x_r; \boldsymbol{\mu}_{\setminus r}^{(t)}, S_{\setminus r}^{(t)}) dx_r \\ \sigma_r^{(t)2} &= \int (x_r - m_r^{(t)})^2 p_r(x_r; \boldsymbol{\mu}_{\setminus r}^{(t)}, S_{\setminus r}^{(t)}) dx_r\end{aligned}$$

$$\text{where } p_r(x_r; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) = \frac{1}{Z_r} \exp \left[c_r(x_r) + \frac{1}{2} (\mathbf{J}_r^T K_{\setminus r}^{-1} \mathbf{J}_r) x_r^2 + (\mathbf{J}_r^T K_{\setminus r}^{-1} \boldsymbol{\mu}_{\setminus r}) x_r \right].$$

3. Update $\mu_r^{(t+1)}$ and $s_r^{(t+1)}$ $r = 1, \dots, L$ as follows

$$\begin{aligned}s_r^{(t+1)} &= \frac{1}{\sigma_r^{(t)2}} + \mathbf{J}_r^{(t)T} K_{\setminus r}^{(t)-1} \mathbf{J}_r^{(t)} \\ \mu_r^{(t+1)} &= \frac{m_r^{(t)}}{\sigma_r^{(t)2}} + \mathbf{J}_r^{(t)T} K_{\setminus r}^{(t)-1} \boldsymbol{\mu}_{\setminus r}^{(t)}.\end{aligned}$$

4. Repeat 2 and 3 until convergence.

(We have no idea if this algorithm works or not.)

5 Summary

We have shown the information geometrical framework of the belief propagation algorithm and that the adaptive TAP equations are derived in the similar manner as the equilibrium conditions of the BP algorithm. Here we defined $q(\mathbf{x})$, $p_0(\mathbf{x}; \boldsymbol{\mu}, S)$, and $p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r})$ in a different way from original BP, but the information geometrical view is almost the same for both cases. We have also shown that the free energy, which corresponds to the Bethe free energy in the case of BP, is the adaptive TAP free energy.

In the case of BP, the solution becomes exact when the graph is tree. In this problem, we do not have the same result, but adaptive TAP becomes exact when the target distribution $q(\mathbf{x})$ is a normal distribution. The information geometrical framework immediately shows this since when $q(\mathbf{x})$ is a normal distribution, $q(\mathbf{x}), p_0(\mathbf{x}; \boldsymbol{\mu}, S), p_r(\mathbf{x}; \boldsymbol{\mu}_{\setminus r}, S_{\setminus r}) \in M_0$, and every distribution becomes equivalent.

This memo shows that the information geometrical frameworks of BP and adaptive TAP are equivalent. From this fact, there is a possibility that we can reuse a lot of results developed for BP and inference problem of loopy graphs. One important direction is the algorithm. We can derive the BP and CCCP algorithms to compute the equilibrium of adaptive TAP equations.

References

- [1] Shiro Ikeda, Toshiyuki Tanaka, and Shun-ichi Amari. Information geometry of turbo and low-density parity-check codes. *IEEE Transactions on Information Theory*, 50(6):1097–1114, June 2004.
- [2] Shiro Ikeda, Toshiyuki Tanaka, and Shun-ichi Amari. Stochastic reasoning, free energy, and information geometry. *Neural Computation*, 16(9):1779–1810, September 2004.
- [3] Manfred Oppen and Ole Winther. Adaptive and self-averaging Thouless-Anderson-Palmer mean field theory for probabilistic modeling. *Physical Review E*, 64(5):056131, 2001.